

Inhomogeneities in the universe

Francesco Sylos Labini

Centro Enrico Fermi, Piazza del Viminale 1, 00184 Rome, Italy and
Istituto dei Sistemi Complessi CNR, Via dei Taurini 19, 00185 Rome, Italy

E-mail: sylos@centrofermi.it

Abstract. Galaxy structures represent one of the most challenging observations in cosmology. Whether or not these structures are found to be compatible with the standard model of galaxy formation crucially depends on the a-priori and assumptions encoded in the statistical methods employed to characterize the data and on the a-posteriori hypotheses made to interpret the results. We present strategies to test the most common assumptions, i.e. spatial homogeneity and statistical homogeneity and isotropy. These tests provide evidences that, in the available samples, galaxy distribution is spatially inhomogeneous but statistically homogeneous and isotropic. Different conclusions are obtained through statistical tools based on a-priori assumptions that not verified in the data when directly tested. Thus we find that the observed galaxy structures are not compatible with the standard model of galaxy formation, e.g. LCDM, that predicts spatial homogeneity at small scales (i.e., $r < 10$ Mpc/h), structures of relatively limited size (i.e. $r < 100$ Mpc/h) and anti-correlations at large scales (i.e. $r > 150$ Mpc/h). While the observed inhomogeneities pose a fundamental challenge to the standard picture of cosmology they also represent an important opportunity which may open new directions for many cosmological puzzles.

PACS numbers: 98.65.-r,98.65.Dx,98.80.-k

1. Introduction

One of the cornerstones of modern cosmology is obtained by the observations of the three dimensional distribution of galaxies [1, 2]. In recent years there has been a fast growth of data which has allowed a detailed characterization of galaxy structures at low redshift (i.e., $z < 0.3$) and small scales (i.e., $r < 150$ Mpc/h). While many authors (see e.g., [3, 4, 5, 6, 7, 8, 9]) have claimed that the data are compatible with the theoretical expectations, and thus they are currently used to estimate a number of cosmological parameters, there are some critical issues which have not received the due attention (see e.g., [10, 11, 12, 14, 15, 16]). The critical points concern the a-priori assumptions which are usually used, without being tested, in the statistical analysis of the data and the a-posteriori hypotheses that are used to interpret the results. Among the former, there are the assumptions of spatial homogeneity and of translational and rotational invariance (i.e., statistical homogeneity) which are built in the standard estimators of galaxy correlations [17]. While these estimators are certainly the correct ones to use when these properties are verified, it is an open question of whether the data satisfy these assumptions. It is indeed well known that galaxies are organized into large scale structures, like clusters, filaments and voids, with large fluctuations [18, 19, 20, 21, 22] and it is not obvious a-priori that spatial or statistical homogeneity are satisfied in a sample of any size. Indeed the critical point is that the assumptions of spatial and statistical homogeneity are encoded in the statistical methods and thus these properties are supposed to be satisfied inside (i.e., not at larger scales) any considered sample where these methods are employed [17].

While it is possible that at large enough scales galaxy distribution becomes spatially uniform, this is certainly not the case at small enough scales. Thus it comes the question of the determination of the scale beyond which spatial homogeneity is satisfied. To investigate this problem it is necessary to use statistical methods which are able to directly test spatial homogeneity, and that thus do not assume it a-priori. This is precisely the reason to introduce statistical methods which are more general than the usual ones [17]. However when facing the question of testing spatial homogeneity, one has to consider a number of subtle possibilities which may affect the determinations of statistical quantities via volume averages [11]. In brief the central problem is the stability of finite sample determinations: when a statistical quantity depends on the finite size of the sample then this is not a meaningful and useful estimator of an ensemble average property. A critical analysis of finite-sample volume averages is thus necessary to identify the subtle effects induced by spatial inhomogeneities.

As mentioned above, there is then a second kind of ad-hoc hypotheses which are often used in the interpretation of the statistical analysis. There are some results which are a-priori unexpected but that, a-posteriori, are interpreted as due to some intervening effects in the data, as for example galaxy evolution or selection effects. While it is very well possible that such effects are present in the data, it is however necessary not to parameterize their influence simply through the introduction of ad-hoc functions and

parameters (see e.g. [23, 24, 4] and discussion in [11]), rather to develop focussed tests to understand whether the additional hypotheses are supported by the actual data.

To frame the problem of the comparison of theoretical predictions with data we focus the attention on the real space correlation properties of standard models. The Friedmann-Robertson-Walker (FRW) geometry is derived under the assumptions that matter distribution is exactly translational and rotational invariant [25]. This implies that the matter density is assumed to be constant in a spatial hyper-surface. The FRW metric describes the geometry of the universe in terms of a single function, the scale factor, which obeys to the Friedmann equations [25]. On the top of the mean field one can consider statistically homogeneous and isotropic small-amplitude fluctuations. These furnish the seeds of gravitational clustering which eventually give rise to the structures we observe in the present universe. The growth of fluctuations into non-linear structures is considered to have a negligible effect on the space-time dynamics, which is instead driven by the uniform mean field [26]. The statistical properties of matter density fluctuations have to satisfy an important condition in order to be compatible with the FRW geometry [27, 28]. In its essence, the condition is that fluctuations in the gravitational potential induced by density fluctuations do not diverge at large scales [29, 17]. This situation requires that the matter density field fluctuations decay in the fastest possible way with scale [30]. Correspondingly the two-point correlation function becomes negative at larger scales (i.e., $r > 150$) Mpc/h which imply the absence of larger structures of tiny density fluctuations. Are the large scale structures and fluctuations compatible with such a scenario [31] ?

This paper is organized as follows. In Sect.2 we briefly review the main properties of both spatially homogeneous and inhomogeneous stochastic density fields. All definitions are given in the ensemble sense, and for ergodic processes, in the infinite volume limit. The main features of density fields in standard cosmological models are presented in Sect.3, focusing the attention to real space correlation properties. In the case of a real point distribution (Sect.4) the information that can be exacted from the data is thorough a statistical analysis in a finite sample, and hence through the computation of volume averages. We discuss how to set up a strategy to analyze a point distribution in a finite volume, stressing the sequence of steps that should be considered in order to reduce as much as possible the role of a-priori assumptions encoded in the statistical analysis. The analysis of the galaxy data is presented in Sect.5. We discuss that galaxy distribution, at relatively low redshifts (i.e., $z < 0.3$) and small scales (i.e., $r < 150$ Mpc/h) is characterized by large density fluctuations which correspond to large-scale correlations. We emphasize that by using the standard statistical tools one reaches a different conclusion. This occurs because these methods are based on several important assumptions: some of them, when directly tested are not verified, while others are very strong ad-hoc hypotheses which require a detailed investigation. Finally in Sect.6 we draw our main conclusions.

2. A brief review of the main statistical properties

In this section we review the main probabilistic properties of mass density fields. This means that we consider ensemble averages or, for ergodic cases, volume averages in the infinite volume limit. The probabilistic properties of a distribution are useful to be studied in view of its statistical characterization in finite samples (see Sec.4).

A mass density field can be represented as stationary stochastic process that consists in extracting the value of the microscopic density function $\rho(\vec{r})$ † at any point of the space. This is completely characterized by its probability density functional $\mathcal{P}[\rho(\vec{r})]$. This functional can be interpreted as the joint probability density function (PDF) of the random variables $\rho(\vec{r})$ at every point \vec{r} . If the functional $\mathcal{P}[\rho(\vec{r})]$ is invariant under spatial translations then the stochastic process is *statistically homogeneous* or translational invariant (stationary) [17]. When $\mathcal{P}[\rho(\vec{r})]$ is also invariant under spatial rotation then the density field is *statistically isotropic* [17].

A crucial assumption usually used, when comparing theoretical prediction to data, is that stochastic fields are required to satisfy spatial *ergodicity*. Let us take a generic observable $\mathcal{F} = \mathcal{F}(\rho(\vec{r}_1), \rho(\vec{r}_2), \dots)$ function of the mass distribution $\rho(\vec{r})$ at different points in space $\vec{r}_1, \vec{r}_2, \dots$. Ergodicity implies that $\langle \mathcal{F} \rangle = \overline{\mathcal{F}} = \lim_{V \rightarrow \infty} \overline{\mathcal{F}}_V$, where the symbol $\langle \dots \rangle$ is for the (ensemble) average over different realizations of the stochastic process, and $\overline{\mathcal{F}}_V = \frac{1}{V} \int_V \mathcal{F} dV$ is the spatial average in a finite volume V [17].

2.1. Spatially homogeneous distributions

The condition of *spatial homogeneity (uniformity)* is satisfied if the ensemble average density of the field $\rho_0 = \langle \rho \rangle$ is strictly positive, i.e. for an ergodic stochastic field,

$$\langle \rho \rangle = \lim_{R \rightarrow \infty} \frac{1}{V(R; \vec{x}_0)} \int_{V(R; \vec{x}_0)} \rho(r) d^3r > 0 \quad \forall \vec{x}_0, \quad (1)$$

where R is the linear size of a volume V with center in \vec{x}_0 . It is necessary to carefully test spatial homogeneity before applying the definitions given in this section to a finite sample distribution (see Sect.4). Indeed, for inhomogeneous distributions the *estimation* of the average density substantially differs from its asymptotic value and thus the sample estimation of ρ_0 is biased by finite size effects. Unbiased tests of spatial homogeneity can be achieved by measuring conditional properties (see below).

A distribution is spatially inhomogeneous up to a scale λ_0 if

$$\left| \frac{1}{V(R; \vec{x}_0)} \int_{V(R; \vec{x}_0)} d^3x \rho(\vec{x}) - \rho_0 \right| < \rho_0 \quad \forall R > \lambda_0, \quad \forall \vec{x}_0. \quad (2)$$

This equation defines the homogeneity scale λ_0 which separates the strongly fluctuating regime $r < \lambda_0$ from the regime where fluctuations have small amplitude relative to the asymptotic average.

† We use the symbol $\rho(r)$ for the microscopic mass density and $n(r)$ for the microscopic number density. However in the following sections we consider only the number density, as it is usually done in studies of galaxy distributions. In that case we can simply replace the symbol $\rho(r)$ with $n(r)$ and all the definitions given in this section remain unchanged.

The quantity $\langle \rho(\vec{r}_1)\rho(\vec{r}_2) \rangle dV_1 dV_2$ gives, in a single realization of a stochastic process, the a-priori probability to find two particles simultaneously placed in the infinitesimal volumes dV_1, dV_2 respectively around \vec{r}_1, \vec{r}_2 . The quantity

$$\langle \rho(r_{12}) \rangle_p dV_1 dV_2 = \frac{\langle \rho(\vec{r}_1)\rho(\vec{r}_2) \rangle}{\rho_0} dV_1 dV_2 \quad (3)$$

gives the a-priori probability of finding two particles placed in the infinitesimal volumes dV_1, dV_2 around \vec{r}_1 and \vec{r}_2 with the condition that the origin of the coordinates is occupied by a particle (Eq.3 is the ratio of unconditional quantities, and thus, for the roles of probabilities, it defines a conditional quantity) [17].

For a stationary and spatially homogeneous distribution (i.e., $\rho_0 > 0$), we may define the reduced two-point correlation function as

$$\xi(r_{12}) = \frac{\langle \rho(r_{12}) \rangle_p}{\rho_0} - 1 = \frac{\langle \rho(r_{12}) \rangle}{\rho_0^2} - 1. \quad (4)$$

This function characterizes correlation properties of density fluctuations with small amplitude with respect to the (ensemble) average density.

Let us now discuss the information that can be extracted from $\xi(r)$, when spatial homogeneity has been already proved. Suppose that correlations have a finite range so that in this case we have

$$\xi(r) = A \exp(-r/r_c), \quad (5)$$

where r_c is the *correlation length* of the distribution and A is a constant. Structures of fluctuations have a size determined by r_c . This length scale depends only on the rate of decay of the correlation function. Another characteristic length scale can be defined as $\xi(r_0) = 1$; from Eq.5 we find

$$r_0 = r_c \cdot \log(A), \quad (6)$$

i.e., it depends on the amplitude A of the correlation function §. The two scales r_0 and r_c have a completely distinct meaning: the former marks the crossover from large to small fluctuations while the latter quantifies the typical size of clusters of *small amplitude* fluctuations. When $\xi(r)$ is a power-law function of separation (i.e. $\xi(r) \sim r^{-\gamma}$ with $0 < \gamma < 3$) then the correlation length r_c is infinite and there are clusters of all sizes [17]. (In cosmology the term correlation length is very often used to mean the scale r_0 instead of r_c . This is in our opinion a confusing terminology, and we use the standard one in statistical physics defined in Eqs.5-6.)

The two-point correlation function defined by Eq.4 is simply related to the normalized mass variance in a volume $V(R)$ of linear size R [17]

$$\sigma^2(R) = \frac{\langle M(R)^2 \rangle - \langle M(R) \rangle^2}{\langle M(R) \rangle^2} = \frac{1}{V^2(R)} \int_{V(R)} d^3 r_1 \int_{V(R)} d^3 r_2 \xi(r_{12}). \quad (7)$$

§ When $\xi(r)$ has a more complex behavior than Eq.5, the scale r_0 is different from Eq.6 but simple to computed from its definition, i.e. $\xi(r_0) = 1$. However the finite range of positive correlations generally corresponds to an exponential decay of the type of Eq.5.

The scale r_* at which fluctuations are of the order of the mean, i.e. $\sigma(r_*) = 1$, is proportional to the scale r_0 at which $\xi(r_0) = 1$ and to the scale λ_0 defined in Eq.2.

For spatially uniform systems, when the volume V in Eq.7 is a real space sphere^{||}, it is possible to proceed to the following classification for the scaling behavior of the normalized mass variance at large enough scales [29, 17]:

$$\sigma^2(R) \sim \begin{cases} R^{-(3+n)} & \text{for } -3 < n < 1 \\ R^{-(3+1)} \log R & \text{for } n = 1 \\ R^{-(3+1)} & \text{for } n > 1 \end{cases} . \quad (8)$$

For $-3 < n < 0$ (which corresponds to $\xi(r) \sim r^{-\gamma}$ with $0 < \gamma = 3 + n < 3$), mass fluctuations are “*super-Poisson*”, typical of systems at the critical point of a second order phase transition [17]: there are long-range correlations and the correlation length r_c is infinite. For $n = 0$ fluctuations are Poisson-like and the system is called *substantially Poisson*: there are no correlations (i.e., a purely Poisson distribution) or correlations limited to small scales of the type described by Eq.5, i.e. a finite correlation length. This behavior is typical of many common physical systems, e.g., a homogeneous gas at thermodynamic equilibrium at sufficiently high temperature. Finally for $n \geq 1$ fluctuations are “*sub-Poisson*” or *super-homogeneous* [29, 17] (or hyper-uniform [30]). In this case $\sigma^2(R)$ presents the fastest possible decay for discrete or continuous distributions [29] and the two-point correlation function has to satisfy a global constraint (see Sect.3). Examples are provided, for instance, by the one component plasma, a well-known system in statistical physics [32], and by a randomly shuffled lattice of particles [17, 33].

Note that any *uniform* stochastic process has to satisfy the following condition

$$\lim_{R \rightarrow \infty} \sigma^2(R) = \lim_{R \rightarrow \infty} \frac{1}{V^2(R)} \int_{V(R)} d^3 r_1 \int_{V(R)} d^3 r_2 \xi(r_{12}) = 0 \quad (9)$$

which implies that the average density ρ_0 , in the infinite volume limit, is a well defined concept, i.e. $\rho_0 > 0$ [17].

2.2. Spatially inhomogeneous distributions

A distribution is spatially inhomogeneous in the ensemble (or in the infinite volume limit) sense if $\lambda_0 \rightarrow \infty$. For statistically homogeneous distributions, from Eq.2, we get that the ensemble average density is $\rho_0 = 0$. Thus unconditional properties are not well defined: if we randomly take a finite volume in an infinite inhomogeneous distribution, it typically contains no points. Therefore only conditional properties are well defined, as for instance the average conditional density defined in Eq.3.

For fractal object the average conditional mass included in a spherical volume grows as $\langle M(r) \rangle_p \sim r^D$: for $D < 3$, the average conditional density presents a scaling behavior of the type

$$\langle \rho(r) \rangle_p = \frac{\langle M(r) \rangle_p}{V(r)} \sim r^{D-3} , \quad (10)$$

^{||} The case in which the volume is a Gaussian sphere can be misleading, see e.g. [29]

so that $\lim_{r \rightarrow \infty} \langle \rho(r) \rangle_p = 0$. The hypotheses underlying the derivation of the Central Limit Theorem are violated by the long-range character of spatial correlations, resulting in a PDF of fluctuations that does not follow the Gaussian function [17, 34]. Actually it typically displays “long tails” [35] which can be associated with the divergence of some moments of the distribution.

It is possible to introduce more complex inhomogeneous distributions than Eq.10, for instance the multi-fractal distributions for which the scaling properties are not described by a single exponent, but they change in different spatial locations being characterized by a spectrum of exponents [17]. Another simple (and different !) example is given by a distribution in which the scaling exponent in Eq.10 depends on distance, i.e. $D = D(r) < 3$.

3. Statistical properties of the standard model

As discussed in the introduction, an important constraint must be valid for any kind of initial matter density fluctuation field in the framework of FRW models. This is represented by the condition of super-homogeneity, corresponding in cosmology to the so-called condition of “scale-invariance” of the primordial fluctuations power spectrum (PS)¶ [29]. To avoid confusion, note that in statistical physics the term “scale invariance” is used to describe the class of distributions which are invariant with respect to scale transformations. For instance, a magnetic system at the critical point of transition between the paramagnetic and ferromagnetic phase, shows a two-point correlation function which decays as a non-integrable power law, i.e. $\xi(r) \sim r^{-\gamma}$ with $0 < \gamma < 3$ (super-Poisson distribution in Eq.8). The meaning of “scale-invariance” in the cosmological context is therefore completely different, referring to the property that the mass variance at the horizon scale be constant (see below) [29].

3.1. Basic Properties

Matter distribution in cosmology is assumed to be a realization of a *stationary* stochastic point process that is also spatially uniform. In the early universe the homogeneity scale λ_0 is of the order of the inter-particle distance, and thus negligible, while it grows during the process of structure formation driven by gravitational clustering. The main property of primordial density fields in the early universe is that they are super-homogeneous, satisfying Eq.8 with $n = 1$. This latter property was firstly hypothesized in the seventies [27, 28] and it subsequently gained in importance with the advent of inflationary models in the eighties [29].

In order to discuss this property, let us recall that the initial fluctuations are taken to have Gaussian statistics and a certain PS. Since fluctuations are Gaussian, the knowledge of the PS gives a complete statistical description of the fluctuation field. In a FRW

¶ The PS of density fluctuations is $P(\vec{k}) = \langle |\delta_\rho(\vec{k})|^2 \rangle$, where $\delta_\rho(\vec{k})$ is the Fourier Transform of the normalized fluctuation field $(\rho(\vec{r}) - \rho_0)/\rho_0$ [29].

cosmology there is a fundamental characteristic length scale, the horizon scale $R_H(t)$ that is simply the distance light can travel from the Big Bang singularity $t = 0$ until any given time t in the evolution of the Universe, and it grows linearly with time. Harrison [27] and Zeldovich [28] introduced the criterion that matter fluctuations have to satisfy on large enough scales. This is named the Harrison-Zeldovich criterion (H-Z); it can be written as

$$\sigma^2(R = R_H(t)) = \text{constant}. \quad (11)$$

This condition states that the mass variance at the horizon scale is constant: it can be expressed more conveniently in terms of the PS for which Eq.11 is equivalent to assume $P(k) \sim k$ (the H-Z PS) and that in a spatial hyper-surface $\sigma^2(R) \sim R^{-4}$ [29, 17].

3.2. Physical implications of super-homogeneity

In order to illustrate the physical implications of the H-Z condition, one may consider the gravitational potential fluctuations $\delta\phi(\vec{r})$, which are linked to the density fluctuations $\delta\rho(\vec{r})$ via the gravitational Poisson equation: $\nabla^2\delta\phi(\vec{r}) = 4\pi G\delta\rho(\vec{r})$. From this, transformed to Fourier space, it follows that the PS of the potential $P_\phi(k) = \langle |\delta\hat{\phi}(\vec{k})|^2 \rangle$ is related to the density PS $P(k)$ through the equation $P_\phi(k) \sim \frac{P(k)}{k^4}$. The H-Z condition, $P(k) \sim k$, corresponds therefore to $P_\phi(k) \propto k^{-3}$, so that the variance of the gravitational potential fluctuations, $\sigma_\phi^2(R) \approx \frac{1}{2}P_\phi(k)k^3|_{k=R^{-1}}$, is constant with k [29].

The H-Z condition is a *consistency constraint* in the framework of FRW cosmology. Indeed, the FRW is a cosmological solution for a perfectly homogeneous Universe, about which fluctuations represent inhomogeneous perturbations. If density fluctuations obey to a different condition than Eq.11, and thus $n < 1$ in Eq.8, then the FRW description *will always break down* in the past or future, as the amplitude of the perturbations become arbitrarily large or small. Thus the super-homogeneous nature of primordial density field is a fundamental property independently on the nature of dark matter. This is a very strong condition to impose, and it excludes even Poisson processes ($n = 0$ in Eq.8) [29] for which the fluctuations in gravitational potential diverge at large scales.

3.3. The small scales behavior

Various models of primordial density fields differ for the behavior of the PS at large wavelengths depending on the specific properties hypothesized of the dark matter component. For example, in the Cold Dark Matter (CDM) scenario, where elementary non-baryonic dark matter particles have a small velocity dispersion, the PS decays as a power law $P(k) \sim k^{-2}$ at large k . For Hot Dark Matter (HDM) models, where the velocity dispersion is large, the PS presents an exponential decay at large k . However at small k they both exhibit the H-Z tail $P(k) \sim k$ which is indeed the common feature of all density fluctuations compatible with FRW models. The scale $r_c \approx k_c^{-1}$ at which the PS shows the turnover from the linear to the decaying behavior is fixed to be the size of the horizon at the time of equality between matter and radiation [41].

3.4. The two-point correlation function and super-homogeneity

The super-homogeneity (or H-Z) condition corresponds to the following limit condition

$$\int_0^\infty d^3r \xi(r) = 0, \quad (12)$$

which is another way to reformulate the condition that $\lim_{k \rightarrow 0} P(k) = 0$. This means that there is a fine tuned balance between small-scale positive correlations and large-scale negative anti-correlations [29, 17]. Note that Eq.12 is different, and much stronger, from the requirement that any *uniform* stochastic process has to satisfy, i.e. Eq.9 [17]. In terms of correlation function $\xi(r)$ CDM/HDM models present the following behavior: it is positive at small scales (decaying as $\xi(r) \sim r^{-1}$ for CDM and being almost flat for HDM), it crosses zero at r_c and then it is negative approaching zero with a tail which goes as $-r^{-4}$ (in the region corresponding to $P(k) \sim k$) [17].

3.5. Baryonic acoustic oscillations

To conclude let us mention the baryon acoustic oscillations (BAO) scale [36]. The physical description which gives rise to these oscillations is based on fluid mechanics and gravity: when the temperature of the plasma was hotter than $\sim 10^3$ K, photons were hot enough to ionize hydrogen so that baryons and photons can be described as a single fluid. Gravity attracts and compresses this fluid into the potential wells associated with the local density fluctuations. Photon pressure resists this compression and sets up acoustic oscillations in the fluid. Regions that have reached maximal compression by recombination become hotter and hence are now visible as local positive anisotropies in the cosmic microwave background radiation (CMBR), if the different k -modes are assumed to have the same phase.

For our discussion, the principal point to note is that while k -oscillations are de-localized, in real space the correlation function shows a characteristic feature: $\xi(r)$ has a localized feature at the scale r_{bao} corresponding to the frequency of oscillations in k space. This simply reflects that the Fourier Transform of a regularly oscillating function is a localized function. Formally the scale r_{bao} corresponds to a scale where a derivative of $\xi(r)$ is not continuous [17, 37].

3.6. Size of structures and characteristic scales

There are thus three characteristic scales in the LCDM-type models (see Fig.1). The first is the homogeneity scale which depends on time $\lambda_0 = \lambda_0(t)$, the second is the scale r_c where $\xi(r_c) = 0$ (that roughly corresponds to the scale defined in Eq.5) which is fixed by the initial properties of the matter density field as well as the third scale r_{bao} . As long as the homogeneity scale is smaller than r_{bao}, r_c , these two scales are substantially unchanged by gravitational dynamics: at those scales this is in the linear regime as fluctuations have a small enough amplitude, and it linearly amplifies the initial fluctuations spectrum. The rate of growth of the homogeneity scale can be simply

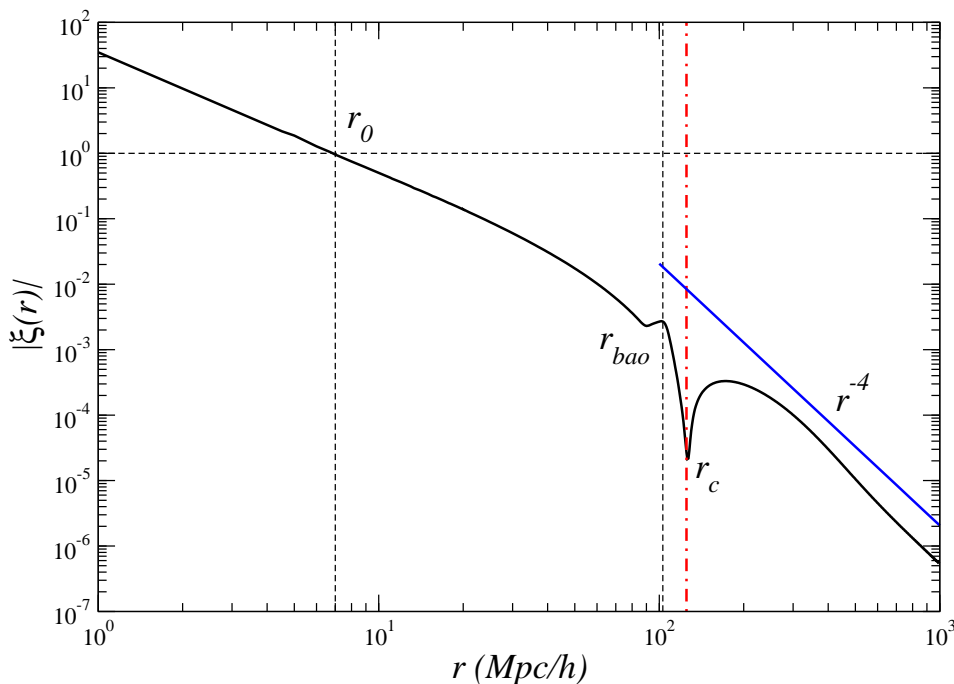


Figure 1. Schematic behavior of the two-point correlation function for the LCDM case. At small scales $r < r_0 \approx 10$ Mpc/h (where $\xi(r_0) = 1$) non-linear gravitational clustering has changed the initial shape of $\xi(r)$. At larger scales $\xi(r)$ has been only amplified by gravitational clustering in the linear regime. For $10 < r < r_c \approx 120$ Mpc/h the correlation is positive and with small amplitude. At larger scales it is negative and characterized by the $\xi(r) \sim -r^{-4}$ behavior. The location of r_{bao} is fixed by cosmological parameters: in the example shown $r_{bao} < r_c$ as predicted by the “concordance model” [3].

computed by using the linear perturbation analysis of a self-gravitating fluid in an expanding universe [26]. Given the initial amplitude of fluctuations and the assumed initial PS of matter density fluctuations, it results that $\lambda_0(t_{now}) \approx 10$ Mpc/h [38].

From the characterization of the two-point correlation function of galaxy distribution we can identify three fundamental tests of standard models ⁺:

- If the homogeneity scale λ_0 is much larger (i.e., a factor 5-10) than ~ 10 Mpc/h, then there is not enough time to form non-linear structures in LCDM models [11].
- If the correlation length r_c (i.e., the zero point of $\xi(r)$) is much larger than ~ 100 Mpc/h then there is a problem in the description of the early universe physics.
- The clear test of inflationary models is given by the detection of the negative part of the correlation function, i.e. the range of scales it behaves as $\xi(r) \sim -r^{-4}$: all models necessarily predict such a behavior ^{*}.

⁺ For the power-spectrum there are additional complications, related how galaxies are biased with respect to the underlying density field: see [39, 31, 40] for further details.

^{*} In the same range of scales the PS is expected to be linear with the wave-number, i.e. $P(k) \sim k$. However selection effects may change the behavior of the PS to constant but not the functional behavior

4. Testing assumptions in the statistical methods

A number of different statistics, determined by making a volume average in a finite sample, can be used to characterize a given distribution. In addition, each statistical quantity can be measured by using different estimators. For this reason we have to set up a strategy to attack the problem if we do not know a-priori which are the properties of the given finite sample distribution. To approach the problem we have to reduce as much as possible the number of a-priori assumptions used the statistical methods, to get the correct information from the data.

We limit our discussion to the case of interest, i.e. a set of N point particles (i.e. galaxies) in a volume V . The microscopic number density can be simply written as $n(\vec{r}) = \sum_i^N \delta^3(\vec{r} - \vec{r}_i)$, where $\delta^3(\vec{r})$ is the Dirac delta function. The statistical quantities defined in Sect.2 can be rewritten in terms of the stochastic variable

$$N_i(V) = \int_{V(\vec{y}_i)} d^3x n(\vec{x}), \quad (13)$$

where \vec{y}_i identifies the coordinates of the center of the volume V . If the center \vec{y}_i coincides with a point particle position \vec{r}_i , then Eq.13 is a conditional quantity. Instead, if the center \vec{y}_i can be any point of space (occupied or not by a particle) then the statistics in Eq.13 is unconditional and it is useful to compute, for instance, the mass variance defined in Eq.7.

For inhomogeneous distributions, unconditional properties are ill-defined (Sect.2) and thus we firstly analyze conditional quantities to then pass, only when in which spatial homogeneity has been detected inside the given sample, to unconditional ones. Therefore, in what follows we take in Eq.13 as volume V a sphere of radius r with center in a distribution point particle, i.e. we consider the stochastic variable defined by the number of points in a sphere \sharp of radius r centered on the i^{th} point of the given set, i.e. $V = V(r; \vec{r}_i)$. The PDF $P(N(r)) = P(N; r)$ of the variable $N_i(r)$ (at fixed r) contains, in principle, information about moments of any order [42]. The first moment is the average conditional density and the second moment is the conditional variance [11].

However before considering the moments of the PDF we should study whether they represent statistically meaningful estimates. Indeed, in the determination of statistical properties through volume averages, one implicitly assumes that statistical quantities measured in different regions of the sample are stable, i.e., that fluctuations in different sub-regions are described by the same PDF. Instead, it may occur that measurements in different sub-regions show systematic (i.e., not statistical) differences, which depend, for instance, on the spatial position of the specific sub-regions. In this case the considered statistic is not statistically stationary in space and its whole-sample average value (i.e., any finite-sample estimation of the PDF moments) is not a meaningful descriptor.

of $\xi(r)$ [39, 31, 40].

\sharp When we take a spherical shell instead of a sphere, then we define a differential quantity instead of an integral one.

4.1. Self-averaging

A simple test to determine whether there are systematic finite size effects affecting the statistical analysis in a given sample of linear size L consists in studying the PDF of $N_i(r)$ in sub-samples of linear size $\ell < L$ placed in different spatial regions of the sample (i.e., S_1, S_2, \dots, S_N). When, at a given scale $r < \ell$, $P(N(r), \ell; S_i)$ is the same, modulo statistical fluctuations, in the different sub-samples, i.e.,

$$P(N(r); \ell; S_i) \approx P(N(r); \ell; S_j) \quad \forall i \neq j, \quad (14)$$

it is possible to consider whole sample average quantities. When determinations of $P(N(r); \ell; S_i)$ in different regions S_i show *systematic* differences, then whole sample average quantities are ill defined. In general, this situation may occur because: (i) the lack of the property of translational invariance or (ii) the breaking of the property of self-averaging due to finite-size effects induced by large-scale structures/voids (i.e., long-range correlated fluctuations).

While the breaking of translational invariance imply the lack of self-averaging property the reverse is not true. For instance suppose that the distribution is spherically symmetric, with origin at r_* and characterized by a smooth density profile, function of the distance from r_* [15]. The average density in a certain volume V , depends on the distance of it from r_* : there is thus a systematic effect and Eq.14 is not satisfied. On the other hand when a finite sample distribution is dominated by a single or by a few structures then, even though it is translational invariant in the infinite volume limit, a statistical quantity characterizing its properties in a finite sample can be substantially affected by finite size fluctuations. For instance, a systematic effect is present when the average (conditional) density largely differs when it is measured into two disjointed volumes placed at different distances from the relevant structures (i.e., fluctuations) in the sample. In a finite sample, if structures are large enough, the measurements may differ much more than a statistical scattering ††. That systematic effect sometimes is referred to as cosmic variance [22] but that is more appropriately defined as breaking of self-averaging properties [11], as the concept of variance (which involves already the computation of an average quantity) maybe without statistical meaning in the circumstances described above [11]. In general, in the range of scales in which statistical quantities give sample-dependent results, then they do not represent fair estimations of asymptotic properties of the given distribution [11].

4.2. Spatial homogeneity

The self-averaging test (Eq.14) is the first one to understand whether a distribution is spatially homogeneous or not inside a given sample. As long as the PDF $P(N, r)$ does not satisfy Eq.14 then the distribution is not only spatially inhomogeneous, but the moments of the PDF are not useful estimators of the underlying statistical properties.

††The determination of statistical errors in a finite volume is also biased by finite size effects [31, 16]

Suppose that Eq.14 is found to be satisfied up to given scale $r < L$. Now we can ask the question: is spatial homogeneity reached for $r < L$?

As mentioned in Sect.2, to this aim it is necessary to employ statistical quantities that do not require the assumption of spatial homogeneity such as conditional ones [17, 11]. Particularly the first moment of $P(N, r)$ provides an estimation of the average conditional density defined in Eq.3, which can be simply written as

$$\overline{n(r)_p} = \frac{1}{M(r)} \sum_{i=1}^{M(r)} \frac{N_i(r)}{V(r)} = \frac{1}{M(r)} \sum_{i=1}^{M(r)} n_i(r). \quad (15)$$

We recall that $N_i(r)$ gives the number of points in a sphere of radius r centered on the i^{th} point and the sum is extended to the all $M(r)$ points contained in the sample for which the sphere of radius r is fully enclosed in the sample volume (this quantity is r dependent because of geometrical constraints, see, e.g., [11]). Analogously to Eq.15 the estimator of the conditional variance can be written as

$$\overline{\sigma_p^2(r)} = \frac{1}{M(r)} \sum_{i=1}^{M(r)} n_i^2(r) - \overline{n(r)_p}^2. \quad (16)$$

When, at the scales $< r$, self-averaging properties are satisfied, one may study the scaling properties of $\overline{n(r)_p}$ and of $\overline{\sigma_p^2(r)}$. As long as $\overline{n(r)_p}$ presents a scaling behavior as a function of spatial separation r , as in Eq.10 with $D < 3$, the distribution is spatially inhomogeneous. When $\overline{n(r)_p} \approx \text{const.}$ then this constant provides an estimation of the ensemble average density and the scale λ_0 where the transition to a constant behavior occurs, marks the homogeneity scale. Only in this latter situation it is possible to study the correlation properties of weak amplitude fluctuations. This can be achieved by considering the function $\xi(r)$ defined in Eq.4.

4.3. The two-point correlation function

Before proceeding, let us clarify some general properties of a generic statistical estimator which are particularly relevant for the two-point correlation function $\xi(r)$. As mentioned above, in a finite sample of volume V we are only able to compute a statistical estimator $\overline{X_V}$ of an ensemble average quantity $\langle X \rangle$. The estimator is valid if

$$\lim_{V \rightarrow \infty} \overline{X_V} = \langle X \rangle. \quad (17)$$

If the ensemble average of the finite volume estimator satisfies

$$\langle \overline{X_V} \rangle = \langle X \rangle \quad (18)$$

the estimator is unbiased. When Eq.18 is not satisfied then there is a systematic offset which has to be carefully considered. Note that the violation of Eq.14 implies that Eq.18 is not valid as well. Finally the variance of an estimator is $\sigma_V^X = \langle \overline{X_V^2} \rangle - \langle \overline{X_V} \rangle^2$. The results given by an estimator must be discussed carefully considering its bias and its variance in any finite sample. A strategy to understand what is the effect of these features consists in changing the sample volume V and study finite size effects [17, 31, 11]. This is crucially important for the two-point correlation function $\xi(r)$ as

any estimator $\overline{\xi(r)}$ of it is generally biased, i.e. it does not satisfy Eq.18 [31, 43]. This is because the estimation of the sample mean density is biased when correlations extend over the whole sample size. Indeed, the most common estimator of the average density is

$$\bar{n} = \frac{N}{V}, \quad (19)$$

where N is the number of points in a sample of volume V . It is simple to show that its ensemble average value can be written as [31]

$$\langle \bar{n} \rangle = \langle n \rangle \left(1 + \frac{1}{V} \int_V d^3r \xi(r) \right). \quad (20)$$

Therefore only when $\xi(r) = 0$ (i.e., for a Poisson distribution), Eq.19 is an unbiased estimator of the ensemble average density: otherwise the bias is determined by the integral of the ensemble average correlation function over the volume V .

The most simple estimator of $\xi(r)$ is the Full-Shell (FS) estimator [31] that can be simply written, by following the definition given in Eq.4, as

$$\overline{\xi(r)} = \frac{\overline{(n(r))_p}}{\bar{n}} - 1, \quad (21)$$

where $\overline{(n(r))_p}$ is the estimator of the conditional density in spherical shells rather than in spheres as for the case of Eq.15. Suppose that in a spherical sample of radius R_s , to estimate the sample density, instead of Eq.19, we use the estimator

$$\bar{n} = \frac{3}{4\pi R_s^3} \int_0^{R_s} \overline{(n(r))_p} 4\pi r^2 dr. \quad (22)$$

Then, the estimator defined by Eq.21 must satisfies the following integral constraint

$$\int_0^{R_s} \overline{\xi(r)} r^2 dr = 0. \quad (23)$$

This condition is *satisfied independently of the functional shape of the underlying correlation function* $\xi(r)$. Thus the integral constraint for the FS estimator does not simply introduce an offset, but it causes a change in the shape of $\overline{\xi(r)}$ for $r \rightarrow R_s$. Other choices of the sample density estimator [31, 43] and/or of the correlation function introduce distortions similar to that in Eq.23.

In order to clarify the effect of the integral constraint for the FS estimator, let us rewrite the ensemble average value of the FS estimator (i.e., Eq.21) in terms of the ensemble average two-point correlation function

$$\langle \overline{\xi(r)} \rangle = \frac{1 + \xi(r)}{1 + \frac{3}{R_s^3} \int_0^{R_s} \xi(r) r^2 dr} - 1. \quad (24)$$

By writing Eq.24 we assume that the stochastic noise is negligible, which of course is not a good approximation at any scale. However in this way we may be able to single out the effect of the integral constraint for the FS estimator. From Eq.24 it is clear that this estimator is biased, as it does not satisfy Eq.18 but only Eq.17.

As an illustrative example, let us now consider the case in which the theoretical $\xi(r)$ is a given by LCDM model. The (ensemble average) estimator given by Eq.24,

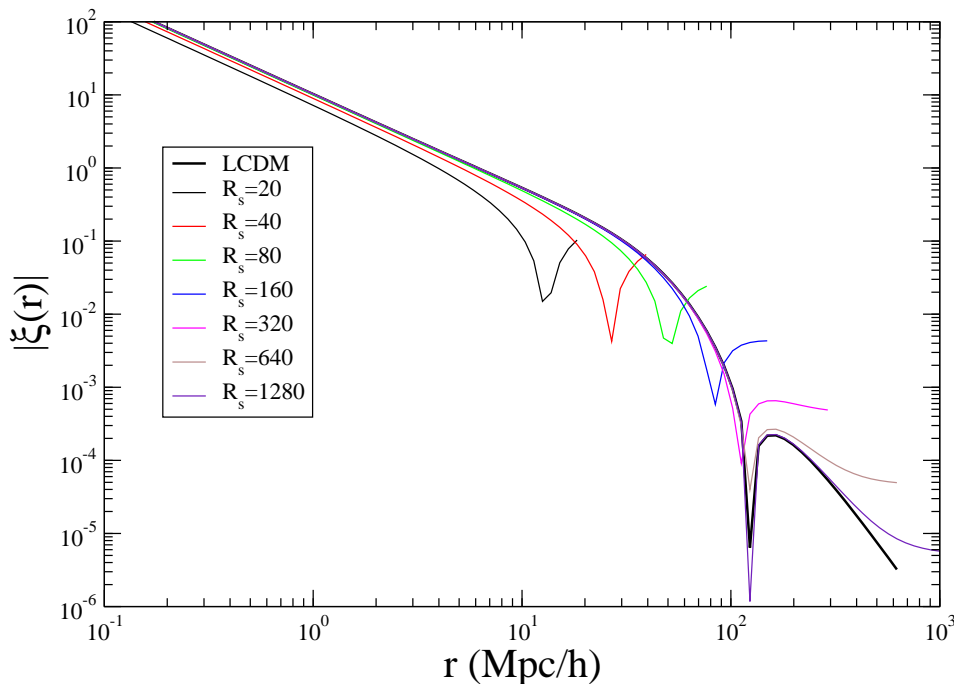


Figure 2. Absolute value of the estimation of the correlation function of the LCDM model with the integral constraint described by Eq.24. The tick solid line represents the theoretical model (From [31]). The zero crossing scale correspond to the cusp.

in spherical samples of different radius R_s , is shown in Fig.2. One may notice that for $R_s > r_c$ the zero point of $\overline{\xi(r)}$ remains stable, while when $R_s < r_c$ it is linearly dependent on R_s . The negative tail continues to be non-linearly distorted even when $R_s > r_c$. For instance, when $R_s \approx 600$ Mpc/h we are not able to detect the $\xi(r) \sim -r^{-4}$ tail that becomes marginally visible only when $R_s > 1000$ Mpc/h. Thus the stability of the zero-point crossing scale should be the first problem to be considered in the analysis of $\overline{\xi(r)}$.

5. Results in the data

We briefly review the main results obtained by analyzing several samples of the Sloan Digital Sky Survey (SDSS) [10, 12, 11, 34, 15] and of the Two degree field Galaxy Redshift Survey (2dFGRS) [44, 13, 14]. For the different catalogs we selected, in the angular coordinates, a sky region such that (i) it does not overlap with the irregular edges of the survey mask and (ii) it covers a contiguous sky area. We computed the metric distance $R(z; \Omega_m, \Omega_\Lambda)$ from the redshift z by using the cosmological parameters $\Omega_m = 0.25$ and $\Omega_\Lambda = 0.75$.

The SDSS catalog includes two different galaxy samples constructing by using different selection criteria: the main-galaxy (MG) sample and the Luminous Red Galaxy (LRG) sample. In particular, the MG sample is a flux limited catalog with apparent magnitude $m_r < 17.77$ [45], while the LRG sample was constructed to be volume-limited

(VL) [46]. A sample is flux limited when it contains all galaxies brighter than a certain apparent flux f_{min} . There is an obvious selection effect in that it contains intrinsically faint objects only when these are located relatively close to the observer, while it contains intrinsically bright galaxies located in wide range of distances [6]. For this reason one constructs a volume limited (VL) sample by imposing a cut in absolute luminosity L_{min} and by computing the corresponding cut in distance $r_{max} \approx \sqrt{L_{min}/(4\pi f_{min})}$, so that all galaxies with $L > L_{min}$, located at distances $r < r_{max}$, have flux $f > f_{min}$, and are thus included in the sample. By choosing different cuts in absolute luminosity one obtains several VL samples (with different L_{min}, r_{max}). Note that we use magnitudes instead of luminosities and that the absolute magnitude must be computed from the redshift by taking into account both the assumptions on the cosmology (i.e. the cosmological parameters, which very weakly perturb the final results given that the redshifts are low, i.e. $z < 0.2$) and the K-corrections (which are measured in the SDSS case).

For the MG sample we used standard K-corrections from the VAGC data [47]: we have tested that our main results do not depend significantly on K-corrections and/or evolutionary corrections [11]. The MG sample angular region we consider is limited, in the SDSS internal angular coordinates, by $-33.5^\circ \leq \eta \leq 36.0^\circ$ and $-48.0^\circ \leq \lambda \leq 51.5^\circ$: the resulting solid angle is $\Omega = 1.85$ sr. For the LRG sample, we exclude redshifts $z > 0.36$ and $z < 0.16$ (where the catalog is known to be not complete [45, 4]), so that the distance limits are: $R_{min} = 465$ Mpc/h and $R_{max} = 1002$ Mpc/h. The limits in R.A α and Dec. δ considered are: $\alpha \in [130^\circ, 240^\circ]$ and $\delta \in [0^\circ, 50^\circ]$. The absolute magnitude is constrained in the range $M \in [-23.2, -21.2]$. With these limits we find $N = 41833$ galaxies covering a solid angle $\Omega = 1.471$ sr [48]. Finally for 2dFGRS, to avoid the effect of the irregular edges of the survey we selected two rectangular regions whose limits are [14]: in southern galactic cap (SGC) ($-33^\circ < \delta < -24^\circ$, $-32^\circ < \alpha < 52^\circ$), and in northern galactic cap (NGC) ($-4^\circ < \delta < 2^\circ$, $150^\circ < \alpha < 210^\circ$); we determined absolute magnitudes M using K-corrections from [49, 14].

5.1. Redshift selection function

In order to have a simple picture of the redshift distribution in a magnitude limited sample, we report Fig.3 galaxy counts as a function of the radial distance, in bins of thickness 10 Mpc/h, in the northern and southern part of the 2dFGRS [14, 13]. One may notice that a sequence of structures and voids is clearly visible, but there is an overall trend (a rise, a peak and then a decrease of the density) which is determined by a luminosity selection effect. Indeed, $n(R)$ in a flux limited sample is usually called redshift selection function, as it is determined by both the redshift distribution and by the luminosity selection criteria of the survey. It is thus not easy, by this kind of analysis, to determine, even at a first approximation, the main properties of the galaxy distributions in the samples. Nevertheless, one may readily compute that there is a $\sim 30\%$ of difference in the sample density between the northern and the southern part of the catalog: one needs to refine the analysis to clarify its significance. Note that

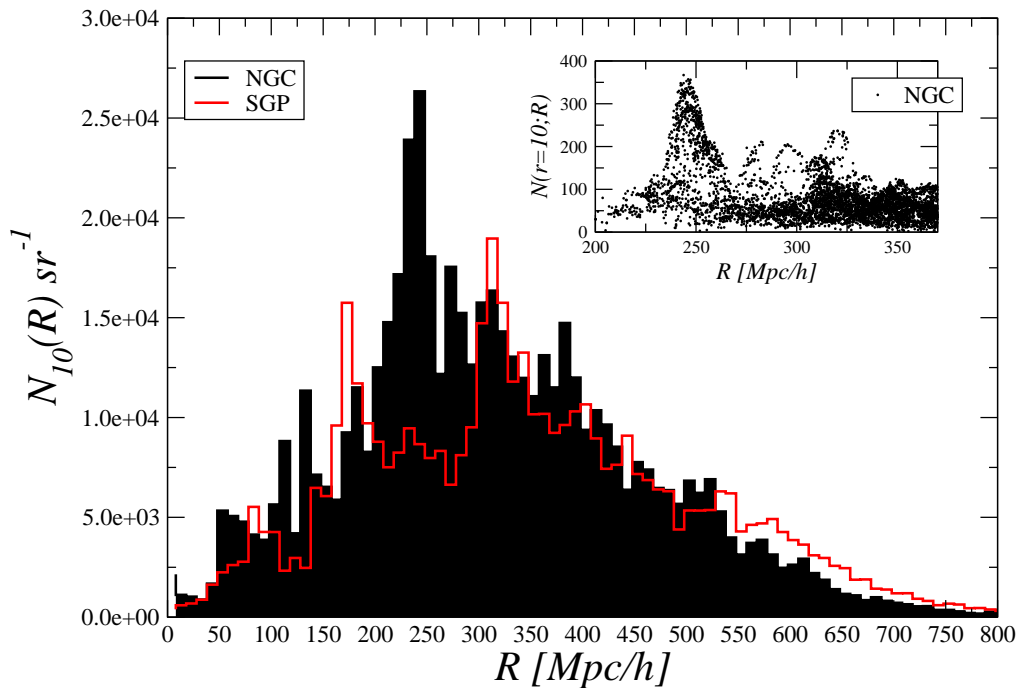


Figure 3. Radial density in bins of thickness 10 Mpc/h in the northern (NGC) and southern (SGP) part of the 2dFGRS magnitude limited sample. There is a large structure at ~ 240 Mpc/h. In the inset panel it is shown the distribution of $N_i(r; R)$ for $r = 10$ Mpc/h in a VL sample in the NGC. (Adapted from [14]).

large scale $\sim 30\%$ fluctuations are not uncommon. For instance, fluctuations have been found in galaxy redshift and magnitude counts that are close to 50% occurring on ~ 100 Mpc/h scales [18, 19, 20, 21].

5.2. Radial counts

A more direct information about the value of the density in a VL sample, is provided by the number counts of galaxies as a function of radial distance $n(R)$ in a VL sample. For a spatially homogeneous distribution $n(R)$ should be constant while, for a fractal distribution it should exhibit a power-law decay, even though large fluctuations are expected to occur given that this not an average quantity [50].

In the SDSS MG VL samples, at small enough scales, $n(R)$ (see Fig.4) shows a fluctuating behavior with peaks corresponding to the main structures in the galaxy distribution [11]. At larger scales $n(R)$ increases by a factor 3 from $R \approx 300$ Mpc/h to $R \approx 600$ Mpc/h. Thus there is no range of scales where one may approximate $n(R)$ with a constant behavior. The open question is whether the growth of $n(R)$ for $R > 300$ Mpc/h is induced by structures or it is caused by a selection effect in data. Both are possible but both must be very detailed discussed. For instance in [24] it is argued that a substantial evolution causes that growth, while in [11] it is discussed, by making a more complete analysis (see next section), that structures certainly contribute to such

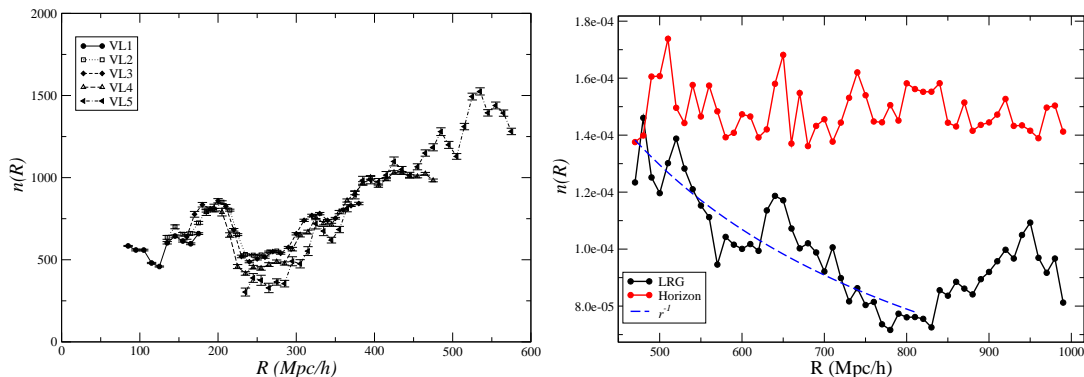


Figure 4. *Left panel:* Radial density in the volume limited samples of the MG catalog. Note the amplitude of $n(R)$ for the MG VL samples has been normalized by taking into account the different selection in luminosity in the different samples (From [11]). *Right Panel:* The same for the LRG sample and for a mock sample extracted from the Horizon simulations [51] (units are in $(\text{Mpc}/h)^{-3}$). The blue dashed line decays as r^{-1} and it is plotted as reference. (From [48]).

a behavior. (Note that in mock catalogs drawn from cosmological N-body simulations one measures an almost constant density [11, 14]).

Given that, by construction, also the LRG sample should be VL [52, 7, 4] the behavior of $n(R)$ is expected to be constant if galaxy distribution is close to uniform (up to Poisson noise and radial clustering). It is instead observed that the LRG sample $n(R)$ shows an irregular and not constant behavior (see the right panel of Fig.4) rather different from that seen in the MG sample. Indeed, there are two main features: (i) a negative slope between $400 \text{ Mpc}/h < r < 800 \text{ Mpc}/h$ (i.e., $0.16 < z < 0.28$) and (ii) a positive slope up to a local peak at $r \sim 950 \text{ Mpc}/h$ (i.e., $z \sim 0.34$). Note that if $n(R)$ were constant we would expect a behavior similar to the one shown by the mock sample extracted from the Horizon simulation [51] with the same geometry of the real LRG sample (see Fig.4) [48].

An explanation that it is usually given for this result [7, 4], is that the LRG sample is “quasi” VL, in that it does not show a constant $n(R)$. Thus, the features in $n(R)$ are absorbed in the properties of a selection function, which is unknown *a priori*, but that it is defined *a posteriori* as the difference between an almost constant $n(R)$ and the behavior observed. This explanation is unsatisfactory as it is given *a posteriori* and no independent tests have been provided to corroborate the hypothesis that an important observational selection effect occurs in the data, other than the behavior of $n(R)$ itself. A different possibility is that the behavior of $n(R)$ is determined, at least partially, by intrinsic fluctuations in the distribution of galaxies and not by selection effects.

By addressing the behavior of $n(R)$ to unknown selection effects, it is implicitly assumed that more than the 20% of the total galaxies have not been measured for observational problems [48]. This looks improbable [52] although a more careful investigation of the problem must be addressed. Note also that the deficit of galaxies

would not be explained by a smooth redshift-dependent effect, rather the selection must be strongly redshift dependent as the behavior of $n(R)$ is not monotonic. These facts point, but do not proof, toward an origin of the $n(R)$ behavior due to the intrinsic fluctuations in the galaxy distribution.

5.3. Test on self-averaging properties

Galaxy counts provide only a rough analysis of fluctuations, especially because one is unable to average it and because it samples different scales differently as the volume in the different redshift bins is not the same. The analysis of the stochastic variable represented by the number of points in spheres $N_i(r)$ can help to overcome these problems, as it is possible to construct volume averages and because it is computed in a simple real sphere. (See an example in the inset panel of Fig.3).

Let us thus pass to the self-averaging test described in Sect.4.1. To this aim we divide the sample into two non-overlapping regions of equal volume, one at low (L) and the other at high (H) redshifts. We then measure the PDF $P_L(N;r)$ and $P_H(N;r)$ in the two volumes. Given that the number of independent points is not very large at large scales (i.e., $M(r)$ in Eq.15 not very larger than $\sim 10^4$), in order to improve the statistics especially for large sphere radii, we allow a partial overlapping between the two sub-samples, so that galaxies in the L (H) sub-sample count also galaxies in the H (L) sub-sample. This overlapping clearly can only smooth out differences between $P_L(N;r)$ and $P_H(N;r)$.

We first consider two SDSS MG VL samples from the data release 6 (DR6) [11] and then from the DR7 [15]. In a first case (upper - left panels of Fig.5), at small scales ($r = 10$ Mpc/h), the distribution is self-averaging (i.e., the PDF is statistically the same) both in the DR6 sample (that covers a solid angle $\Omega_{DR6} = 0.94$ sr) than in the DR7 sample ($\Omega_{DR7} = 1.85$ sr $\approx 2 \times \Omega_{DR6}$ sr). Instead, for larger sphere radii i.e., $r = 80$ Mpc/h, (bottom - right panels of Fig.5) in the DR6 sample, the two PDF show clearly a systematic difference. Not only the peaks do not coincide, but the overall shape of the PDF is not smooth displaying a different shape. Instead, for the sample extracted from DR7, the two determinations of the PDF are in good agreement (within statistical fluctuations). We conclude that in DR6 for $r = 80$ Mpc/h there are large density fluctuations which are not self-averaging because of the limited sample volume [11, 15]. They are instead self-averaging in DR7 because the volume is increased by a factor two.

For the other sample we consider, which include mainly bright galaxies, the breaking of self-averaging properties does not occur as well for small r but it is found for large r . Other radial distance-dependent selections, like galaxy evolution [24], could in principle give an effect in the same direction if they increase the number density with redshift. However this would not affect the conclusion that, on large enough scales, self-averaging is broken. Note that in the SDSS samples for small values of r the PDF is found to be statistically stable in different sub-regions of a given sample. For this reason we do not

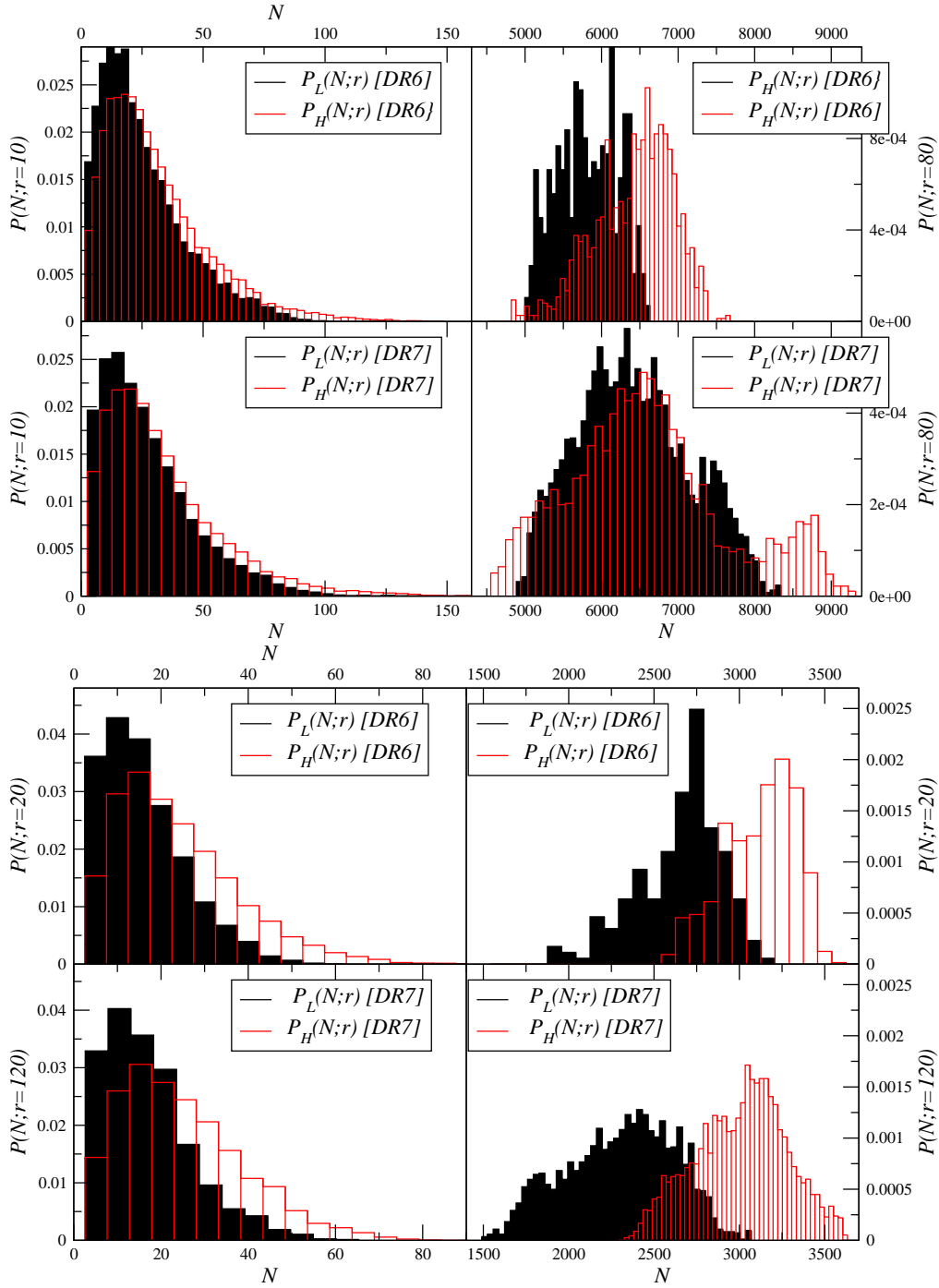


Figure 5. *Upper Panels:* PDF of the counts in spheres in the sample defined by $R \in [125, 400]$ Mpc/h and $M \in [-20.5, -22.2]$ in the DR6 and DR7 data, for two different values of the sphere radii $r = 10$ Mpc/h and $r = 80$ Mpc/h. *Lower Panels:* The same but for the sample defined by $R \in [200, 600]$ Mpc/h and $M \in [-21.6, -22.8]$ and for $r = 20, 120$ Mpc/h. (Adapted from [15]).

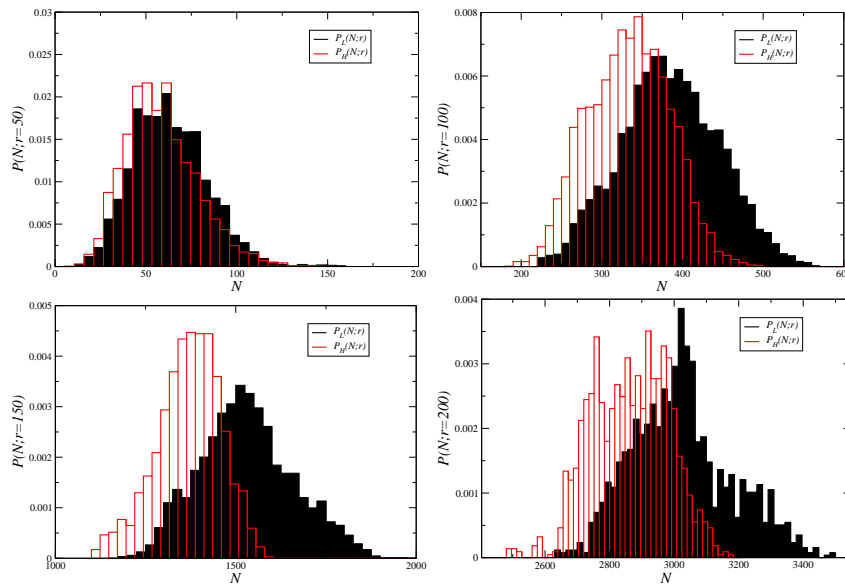


Figure 6. *Upper Left Panel:* PDF for $r = 50$ Mpc/h in the LRG sample (From [48]). The number of points contributing to the histogram is respectively for L and H $M(r) = 13277, 13099$. *Upper Right Panel:* The same for $r = 100$ Mpc/h. Here $M(r) = 7929, 7690$. *Bottom Left Panel:* The same for $r = 150$ Mpc/h. Here $M(r) = 3495, 3150$. *Bottom Right Panel:* The same for $r = 200$ Mpc/h. Here $M(r) = 1465, 1354$.

interpret the lack of self-averaging properties as due to a “local hole” around us: this would affect all samples and all scales, which is indeed not the case [15]. Because of these large fluctuations in the galaxy density field, self-averaging properties are well-defined only in a limited range of scales where it is then statistically meaningful to measure whole-sample average quantities [11, 34, 15].

For the LRG sample (see Fig.6) one may note that for $r = 50$ Mpc/h the two determinations are much closer than for larger sphere radii for which there is actually a noticeable difference in the whole shape of the PDF. The fact that $P_H(N; r)$ is shifted toward smaller values than $P_L(N; r)$ is related to the decaying behavior of the redshift counts (see Fig.4): most of the galaxies at low redshifts see a relatively larger local density than the galaxies at higher redshift.

Due the breaking of self-averaging properties in the different samples for $r < 150$ Mpc/h we conclude that there is no evidence for a crossover to spatial uniformity. In the next section we refine the analysis for smaller scales by characterizing the shape of the PDF and the scaling of its moments.

5.4. Probability density function and its moments

In the range of scales in which self-averaging properties are found to hold, we can further characterize the shape of the PDF and the scaling of its moments. We first computed the average conditional density (Eq.15) finding a pronounced r dependence, as can be seen

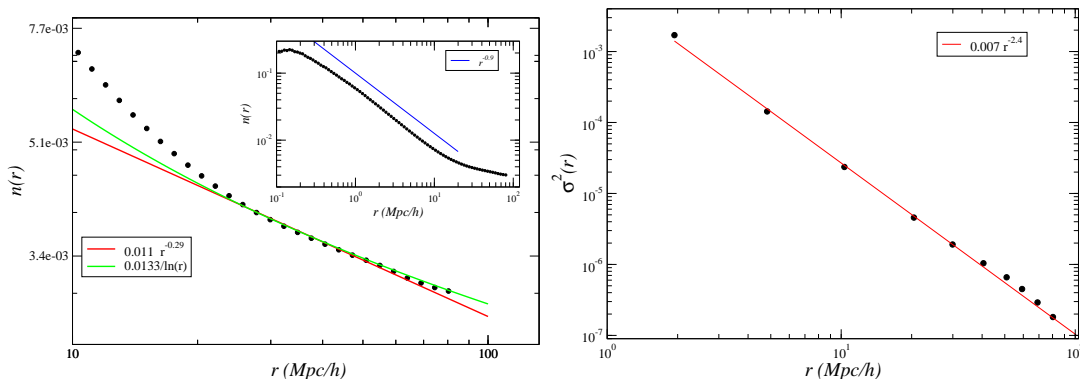


Figure 7. *Left Panel:* Conditional average density $\bar{n}(r)$ of galaxies as a function of radius. In the inset panel it is shown the behavior of in the full range of scales. Note the change of slope at ≈ 20 Mpc/h and note also that there is no flattening up to ≈ 80 Mpc/h. The statistical significance of the last few points at the largest scales is weaker (see text). Our conjecture is that we have a logarithmic correction to the constant behavior, although we cannot exclude the possibility that it is power law with an exponent ≈ -0.3 *Right Panel:* Variance σ^2 of the conditional density $n_i(r)$ as a function of radius. Conversely, the variance for a Poisson point process would display a $1/r^3$ decay.

in Fig.7. We detect a change of slope in the conditional average density in terms of the radius r at about ≈ 20 Mpc/h. At this point the decay of the density changes from an inverse linear decay to a slow logarithmic one. Our best fit is $\bar{n}(r) \approx \frac{0.0133}{\log r}$, that is the average density depends only weakly on r . Alternatively, an almost indistinguishable power-law fit is provided by $\bar{n}(r) \approx 0.011 \times r^{-0.29}$. Moreover, the density $\bar{n}(r)$ does not saturate to a constant up to ~ 80 Mpc/h, i.e., up to the largest scales probed in this sample where self-averaging properties have been tested to hold. Our best fit for the variance is ($\sigma^2(r) \approx 0.007 \times r^{-2.4}$ see right panel of Fig.7). Given the scaling behavior of the conditional density and variance, we conclude that galaxy structures are characterized by non-trivial correlations for scales up to $r \approx 80$ Mpc/h.

To probe the whole distribution of the conditional density $n_i(r)$, we fitted the measured PDF with Gumbel distribution via its two parameters α and β [34]. The Gumbel distribution is one of the three extreme value distribution [53, 54]. It describes the distribution of the largest values of a random variable from a density function with faster than algebraic (say exponential) decay. The Gumbel distribution's PDF is given by

$$P(y) = \frac{1}{\beta} \exp \left[-\frac{y - \alpha}{\beta} - \exp \left(-\frac{y - \alpha}{\beta} \right) \right]. \quad (25)$$

With the scaling variable

$$x = \frac{y - \alpha}{\beta} \quad (26)$$

the density function (Eq.25) simplifies to the parameter-free Gumbel

$$P(x) = e^{-x - e^{-x}}. \quad (27)$$

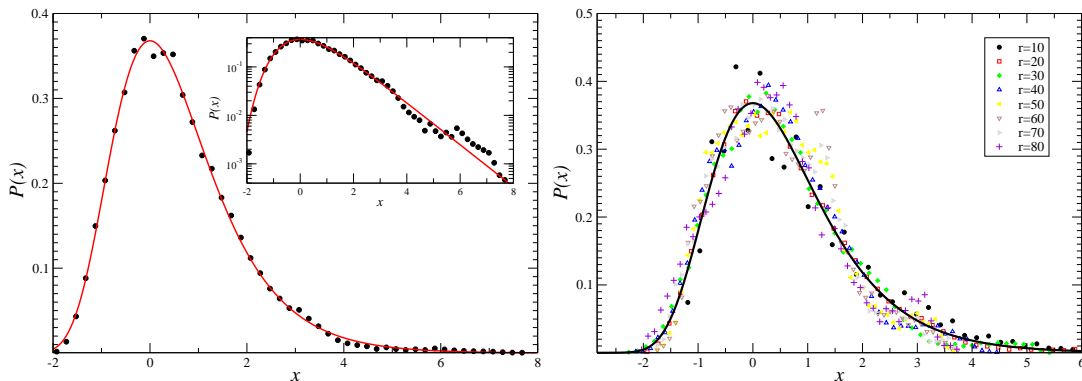


Figure 8. *Left Panel:* One of the best fits is obtained for $r = 20$ Mpc/h. The data is rescaled by the fitted parameters α and β . The solid line corresponds to the parameter-less Gumbel distribution Eq.27. The inset depicts the same on log-linear scale. *Right Panel:* Data curves of different r scaled together by fitting parameters α and β for each curves. The solid line is the parameter-free Gumbel distribution Eq.27.

The mean and the variance of the Gumbel distribution (Eq.25) is $\mu = \alpha + \gamma\beta$, $\sigma^2 = (\beta\pi)^2/6$ where $\gamma = 0.5772\dots$ is the Euler constant.

One of our best fit for the PDF is obtained for $r = 20$ Mpc/h (see left panel Fig. 8). The data, moreover, convincingly collapses to the parameter-less Gumbel distribution (Eq.27) for all values of $r \in [10, 80]$ Mpc/h, with the use of the scaling variable x from Eq.26 (see right panel Fig. 8). Note that for a Poisson point process the number $N(r)$ fluctuations are distributed exactly according to a Poisson distribution, which in turn converges to a Gaussian distribution for large average number of points $\overline{N(r)}$ per sphere. In our samples, $\overline{N(r)}$ was always larger than 20 galaxies, where the Poisson and the Gaussian PDFs differ less than the uncertainty in our data. Note also that due to the Central Limit Theorem, all homogeneous point distributions (not just the Poisson process) lead to Gaussian fluctuations [17]. Hence the appearance of the Gumbel distribution is a clear sign of inhomogeneity and large scale structures in our samples.

We have thus established scaling and data collapse over a wide range of radius (volume) in galaxy data. Scaling in the data indicates criticality [17, 34]. The average galaxy density depends only logarithmically on the radius, which suggests a Gumbel scaling function. Indeed, it was recently conjectured [55] that only three types of distributions appear to describe fluctuations of global observables at criticality. In particular, when the global observable depends logarithmically on the system size, the corresponding distribution should be a (generalized) Gumbel [34].

5.5. Two-point correlation analysis

When one determines the standard two-point correlation function one makes implicitly the assumptions that, inside a given sample the distribution is: (i) self-averaging and (ii) spatially uniform. The first assumption is used when one computes whole sample average quantities. The second is employed when supposing that the estimation of the

sample average gives a fairly good estimation of the ensemble average density. When one of these assumptions, or both, is not verified then the interpretation of the results given by the determinations of the standard two-point correlation function must be reconsidered with great care.

To show how non self-averaging fluctuations inside a given sample bias the $\xi(r)$ analysis, we consider the estimator

$$\overline{\xi(r)} + 1 = \overline{\xi(r; R, \Delta R)} + 1 = \overline{n(r, \Delta r)_p} \cdot \frac{V(r^*)}{N(r^*; R, \Delta R)}, \quad (28)$$

where the second ratio on the r.h.s. is now the density of points in spheres of radius r^* averaged over the galaxies lying in a shell of thickness ΔR around the radial distance R . If the distribution is homogeneous, i.e., $r^* > \lambda_0$, and statistically stationary, Eq.28 should be (statistically) independent on the range of radial distances $(R, \Delta R)$ chosen. The two-point correlation function is defined as a ratio between the average conditional density and the sample average density: if both vary in the same way when the radial distance is changed, then its amplitude remains nearly constant. This however does not imply that the amplitude of $\overline{\xi(r)}$ is meaningful, as it can happen that the density estimated in sub-volumes of size r^* show large fluctuations and so the average conditional density, and this occurring with a radial-distance dependence. To show that the $\overline{\xi(r)}$ analysis gives a meaningful estimate of the amplitude of fluctuations, *one has to test that this amplitude remains stable by changing the relative position of the sub-volumes of size r^* used to estimate the average conditional density and the sample average density.* This is achieved by using the estimator in Eq.28. While standard estimators [56, 43, 31] are not able to test for such an effect, as the main contributions for both the conditional density and the sample average density come from the same part of the sample (typically the far-away part where the volume is larger). We find large variations in the amplitude of $\overline{\xi(r)}$ in the SDSS MG VL samples (see the left panel of Fig.9). This is simply an artifact generated by the large density fluctuations on scales of the order of the sample sizes. The results that the estimator of $\xi(r)$ has nearly the same amplitude in different samples, e.g., [57, 58, 59, 60, 7, 8, 9], despite the large fluctuations of $N_i(r; R)$, are simply explained by the fact that $\overline{\xi(r)}$ is a ratio between the average conditional density and the sample average density: both vary in the same way when the radial distance is changed and thus the amplitude is nearly constant. To summarize the fact that by using different normalizations, which however are all in principle equally valid if the distribution has a well-defined average density inside the sample, we have shown that the amplitude of the estimated correlation function varies in the SDSS samples. This is due to the fact that both the assumptions on which the determination of the standard to point correlation function is based on, are not verified in these samples and that λ_0 is certainly larger than the samples size.

The right panel of Fig.9 clearly show that there is a finite-size dependence of both the amplitude of the correlation function and of the zero-crossing scale: this situations looks like the one shown in Fig.2. Thus prior to the characterization of a fine feature as the BAO scale, it is necessary to test that the correlation function remains stable

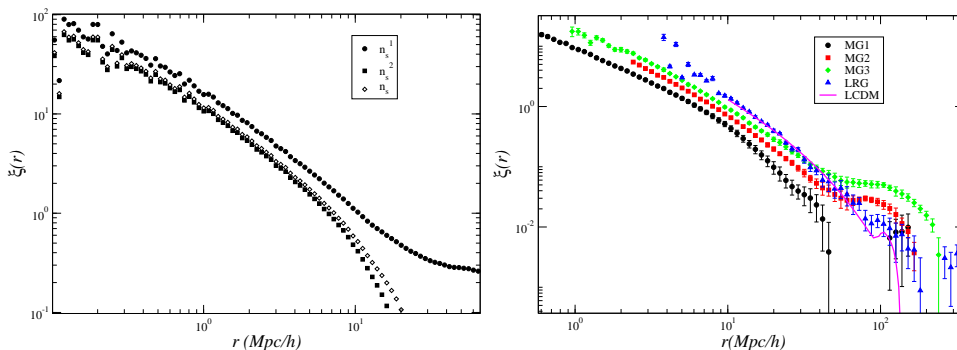


Figure 9. *Left panel:* The two-point correlation function in a MG-VL sample estimated by Eq.28: the sample average density is computed in spheres of radius $r^* = 60$ Mpc/h and considering all center-points lying in a bin of thickness $\Delta R = 50$ Mpc/h centered at different radial distance R : $R_1 = 250$ Mpc/h (n_s^1) and $R_2 = 350$ Mpc/h (n_s^2). The case in which we have used the estimation of the sample average N/V (n_s) is also shown and it agrees with the FS estimator (adapted from [11]). *Right Panel:* The Landy and Szalay [56] estimator of $\xi(r)$ in various MG-VL sample and in a LRG sample of the SDSS. The most evident feature is the finite-size dependence of both the amplitude and the zero-crossing (adapted from [16]). The solid line is a LCDM model.

as a function of the sample size. This shows that estimator of $\xi(r)$ in the LRG sample is indeed biased by volume-dependent systematic effects that make the detection of correlations only an estimate of a lower limit of their amplitude [16]. A similar conclusion was reached by [61], i.e. that when corrections for possible systematics are taken into account the correlation function may not be consistent with as high amplitude a peak as claimed by [3]. To clarify this issue, as discussed above, it is necessary to consider the set of tests for statistical and spatial homogeneity discussed above.

Instead of investigating the origin of the fluctuating behavior of $n(R)$, some authors [4] focused their attention on the effect of the radial counts on the determination of the two-point correlation function $\xi(r)$. In particular, they proposed mainly two different tests to study what is the effect of $n(R)$ on the determination of $\xi(r)$. The first test consists in taking a mock LRG sample, constructed from a cosmological N-body simulation of the LCDM model, and by applying a redshift selection which randomly excludes points in such a way that the resulting distribution has the same $n(R)$ of the real sample. Then one can compare $\xi(r)$ obtained in the original mock and in redshift-sampled mock. [4] find that there is a good agreement between the two. This shows that the particular kind of redshift-dependent random sampling considered for the given distribution, does not alter the determination of the correlation function. Alternatively we may conclude that, under the assumption that the observed LRG sample is a realization of a mock LCDM simulation, the $n(R)$ does not affect the result. However, if we want to test whether the LRG sample has the same statistical properties of the mock catalog, we cannot clearly proof (or disproof) this hypothesis by assuming a priori that this is true.

In other words, standard analyses ask directly the question of whether the data are compatible with a given model, by considering only a few statistical measurements. As it was shown by [5] the LRG correlation function does not pass the null hypothesis, i.e. it are compatible with zero signal, implying that the volume of current galaxy samples is not large enough to claim that the BAO scale is detected. In addition, by *assuming* that the galaxy correlations are modeled by a LCDM model, one may find that the data allow to constrain the position of the BAO scale. In our view this approach is too narrow: in evaluating whether a model is consistent with the data, one should show that *at least the main* statistical properties of the model are indeed consistent with the data. As discussed above, a number of different properties can be considered, which are useful to test the assumptions of (i) self-averaging (ii) spatial homogeneity. When, inside the given sample, the assumption (i) and/or (ii) are/is violated then the compatibility test of the data with a LCDM model is not consistent with the properties of the data themselves.

6. Conclusion

The statistical characterization of galaxy structures presents a number of subtle problems. These are associated both with the a-priori assumptions which are encoded in the statistical methods used in the measurements and in the a-posteriori hypotheses that are invoked to explain certain measured behaviors. By increasing the number of the former, as for example, bias, luminosity evolution, selection effects, one may find that the data are compatible with a certain model. However it is possible to introduce direct tests to understand both whether the a-priori assumptions are compatible with the data and whether it is necessary to introduce a-posteriori untested, but plausible, hypotheses to interpret the results of the data analysis. For instance, the analysis of the simple counts as a function of distance, in the SDSS samples, shows clearly that the observed behavior is incompatible with model predictions. As mentioned above, one may assume that the differences between the model and the observations are due to selection effects. Then this becomes clearly the most important assumption in the data analysis that must be stressed clearly and explicitly. In addition, one may consider whether there is an independent way to study whether there are such strong selection effects in the data.

On the basis of the series of test we have presented, aiming to directly test whether spatial and statistical homogeneity are verified inside the available samples, extracted both from the SDSS and 2dFGRS catalogs, we conclude that galaxy distribution is characterized by structures of large spatial extension. Given that we are unable to find a crossover toward homogeneity, the amplitude of these structures remain undetermined and their main characteristic is represented by the scaling behavior of their relevant statistical properties. In particular, we discussed that the average conditional density presents a scaling behavior of the type $\sim r^{-\gamma}$ with $\gamma \approx -1$ up to ~ 20 Mpc/h followed by a $\gamma \approx -0.3$ behavior up to ~ 100 Mpc/h. Correspondingly the probability density

function (PDF) of galaxy (conditional) counts in spheres show a relatively long tail: it is well fitted by the Gumbel function instead than by the Gaussian function, as it is generally expected for spatially homogeneous density fields.

Our statistical tests can thus provide direct observational tests of the basic assumptions used in the derivation of the FRW models, i.e. spatial and statistical homogeneity. In this respect it is worthing to clarify in more details the subtle difference between these two concepts [15]. A widespread idea in cosmology is that the so-called concordance model of the universe combines *two* fundamental assumptions. The first is that the dynamics of space-time is determined by Einstein’s field equations. The second is that the universe is homogeneous and isotropic. This hypothesis, usually called the Cosmological Principle, is though to be a *generalization* of the Copernican Principle that “the Earth is not in a central, specially favored position” [63, 64]. The FRW model is derived under these two assumptions and it describes the geometry of the universe in terms of a single function, the scale factor, which obeys to the Friedmann equation [25]. There is a subtlety in the relation between the Copernican Principle (all observes are equivalent and there are no special points and directions) and the Cosmological Principle (the universe is homogeneous and isotropic). Indeed, the fact that the universe looks the same, at least in a statistical sense, in all directions and that all observers are alike does not imply spatial homogeneity of matter distribution. It is however this latter condition that allows us to treat, above a certain scale, the density field as a smooth function, a fundamental hypothesis used in the derivation of the FRW metric. Thus there are distributions which satisfy the Copernican Principle and which do not satisfy the Cosmological Principle [17]. These are statistically homogeneous and isotropic distributions which are also spatially inhomogeneous. Therefore the Cosmological Principle represents a specific case, holding for spatially homogeneous distributions, of the Copernican Principle which is, instead, much more general. Statistical and spatial homogeneity refer to two different properties of a given density field. The problem of whether a fluctuations field is compatible with the conditions of the absence of special points and direction can be reformulated in terms of the properties of the PDF which generates the stochastic field.

By analyzing the PDF in the available galaxy samples we can make tests on both the Copernican and Cosmological Principles at low redshift, where we can neglect the important complications of evolving observations onto a spatial surface for which we need a specific cosmological model. We have discussed, however, that the statistical properties of the matter density field up to a few hundreds Mpc is crucially important for the theoretical modeling. We have shown that galaxy distribution in different samples of the SDSS is compatible with the assumptions that this is transitionally invariant, i.e. it satisfies the requirement of the Copernican Principle that there are no spacial points or directions. On the other hand, we found that there are no clear evidences of spatial homogeneity up to scales of the order of the samples sizes, i.e. ~ 100 Mpc/h. This implies that galaxy distribution is not compatible with the stronger assumption of spatial homogeneity, encoded in the Cosmological Principle. In addition, at the largest

scales probed by these samples (i.e., $r \approx 150$ Mpc/h) we found evidences for the breaking of self-averaging properties, i.e. that the distribution is not statistically homogeneous. Forthcoming redshift surveys will allow us to clarify whether on such large scales galaxy distribution is still inhomogeneous but statistically stationary, or whether the evidences for the breaking of spatial translational invariance found in the SDSS samples were due to selection effects in the data.

We note an interesting connection between spatial inhomogeneities and large scale flows which can be hypothesized by assuming that the gravitational fluctuations in the galaxy distribution reflect those in the whole matter distribution, and that peculiar velocities and accelerations are simply correlated. Peculiar velocities provide an important dynamical information as they are related to the large scale matter distribution. By studying their local amplitudes and directions, these velocities allow us, in principle, to probe deeper, or hidden part, of the Universe. The peculiar velocities are indeed directly sensitive to the total matter content, through its gravitational effects, and not only to the luminous matter distribution. However, their direct observation through distance measurements remains a difficult task. Recently, there have been published a growing number of observations of large-scale galaxy coherent motions which are at odds with standard cosmological models [67, 66, 68, 69].

It is possible to consider the PDF of gravitational force fluctuations generated by source field represented by galaxies, and test whether it converges to an asymptotic shape within sample volumes. In several SDSS sample we find that density fluctuations at the largest scales probed, i.e. $r \approx 100$ Mpc/h, still significantly contribute to the amplitude of the gravitational force [65]. Under the hypotheses mentioned above we may conclude that that large-scale fluctuations in the galaxy density field can be the source of the large scale flows recently observed.

As a final remark we mention the growing work to understand the effect of inhomogeneities on the large scale dynamics of the universe [70, 71, 72, 64, 73, 74, 75, 76]. As long as structures are limited to small sizes, and fluctuations have low amplitude, one can just treat fluctuations as small-amplitude perturbations to the leading order FRW approximation. However if structures have “large enough” sizes and “high enough” amplitudes, a perturbation approach may lose its validity and a more general treatment of inhomogeneities needs to be developed. From the theoretical point of view, it is then necessary to understand how to treat inhomogeneities in the framework of General Relativity. To this aim one needs to carefully consider the information that can be obtained from the data. At the moment it is not possible to get some statistical information for large redshifts ($z \approx 1$), but the characterization of relatively small scales properties (i.e., $r < 200$ Mpc/h) is getting more and more accurate. According to FRW models the linearity of Hubble law is a consequence of the homogeneity of the matter distribution. Modern data show a good linear Hubble law even for nearby galaxies ($r < 10$ Mpc/h). This raises the question of why the linear Hubble law is linear at scales where the visible matter is distributed in-homogeneously. Several solution to this apparent paradox have been proposed [72, 77, 78]: this situation shows that already

the small scale properties of galaxy distribution have a lot to say on the theoretical interpretation of their properties. Indeed, while observations of galaxy structures have given an impulse to the search for more general solution of Einstein's equations than the Friedmann one, it is now a fascinating question whether such a more general framework may provide a different explanation to the various effects that, within the standard FRW model, have been *interpreted* as Dark Energy and Dark Matter.

Acknowledgments

I am grateful to T. Antal, Y. Baryshev, A. Gabrielli, M. Joyce M. López-Corredoira and N. L. Vasilyev for fruitful collaborations, comments and discussions. I acknowledge the use of the Sloan Digital Sky Survey data (<http://www.sdss.org>), of the 2dFGRS data (<http://www.mso.anu.edu.au/2dFGRS/>) of the NYU Value-Added Galaxy Catalog (<http://ssds.physics.nyu.edu/>), of the Millennium run semi-analytic galaxy catalog (<http://www.mpa-garching.mpg.de/galform/agnpaper/>)

- [1] York D et al 2000 *Astronom.J.* **120** 1579
- [2] Colless M et al 2001 *Mon.Not.R.Acad.Soc.* **328** 1039
- [3] Eisenstein D J et al 2005 *Astrophys.J.* **633** 560
- [4] Kazin I et al 2010 *Astrophys.J.* **710** 1444
- [5] Cabré A and Gaztanaga E [arXiv:1011.2729](https://arxiv.org/abs/1011.2729)
- [6] Zehavi I et al 2002 *Astrophys.J.* **571** 172
- [7] Zehavi I et al 2005, *Astrophys.J.* **630** 1
- [8] Norberg E et al 2001 *Mon.Not.R.Acad.Soc.* **328** 64
- [9] Norberg E et al 2002 *Mon.Not.R.Acad.Soc.* **332** 827
- [10] Sylos Labini F Vasilyev N L and Baryshev Yu V 2007, *Astron.Astrophys.* **465** 23
- [11] Sylos Labini F Vasilyev N L Baryshev Y V 2009 *Astron.Astrophys* **508** 17
- [12] Sylos Labini F et al 2009 *Europhys.Lett.* **86** 49001
- [13] Sylos Labini F Vasilyev N L Baryshev Y V 2009 *Europhys.Lett.* **85** 29002
- [14] Sylos Labini F Vasilyev N L Baryshev Y V 2009 *Astron.Astrophys.* **496** 7
- [15] Sylos Labini F and Baryshev Y V 2010 *J.Cosmol.Astropart.Phys.* **JCAP06** 021
- [16] Sylos Labini F Vasilyev N L Baryshev Yu V López-Corredoira M 2009 *Astron.Astrophys.* **505** 981
- [17] Gabrielli A Sylos Labini F Joyce M and Pietronero L 2005 *Statistical Physics for Cosmic Structures* (Springer Verlag, Berlin)
- [18] Ratcliffe A Shanks T Parker Q A and Fong R 1998, *Mon.Not.R.Acad.Soc.* **293** 197
- [19] Buswell G S et al 2004 *Mon.Not.R.Acad.Soc.* **354** 991
- [20] Frith W J et al 2003 *Mon.Not.R.Acad.Soc.* **345** 1049
- [21] Frith W J Metcalfe N and Shanks T 2006 *Mon.Not.R.Acad.Soc.* **371** 1601
- [22] Yang A and Saslaw W C [arXiv:1009.0013v1](https://arxiv.org/abs/1009.0013v1)
- [23] Blanton M R et al 2003 *Astrophys.J.* **592** 819
- [24] Loveday J 2004 *Mon.Not.R.Acad.Soc.* **347** 601
- [25] Weinberg S 2008 *Cosmology* (Oxford University Press, Oxford)
- [26] Peebles, P J E *Large Scale Structure of the Universe*, (Princeton University Press, Princeton 1980)
- [27] Harrison E R 1970 *Phys.Rev.D* **1** 2726
- [28] Zeldovich Ya B 1972 *Mon.Not.R.Acad.Soc.* **160** 1
- [29] Gabrielli A Joyce M and Sylos Labini F 2002 *Phys.Rev.* **D65** 083523
- [30] Torquato S and Stillinger F H 2003 *Phys. Rev.* **E 68** 041113
- [31] Sylos Labini F and Vasilyev N L 2008 *Astron.Astrophys* **477** 381
- [32] Gabrielli A et al 2003 *Phys.Rev.* **D67** 043406

- [33] Gabrielli A Joyce M and Torquato S 2008 *Phys.Rev.* **E77** 031125
- [34] Antal T Sylos Labini F Vasilyev N L and Baryshev Y V 2009 *Europhys.Lett.* **88** 59001
- [35] Bouchaud J-P and Potters M 1999 *Theory of Financial Risks* (Cambridge University Press, Cambridge) cond-mat/9905413
- [36] Eisenstein D J and Hu W 1998 *Astrophys.J.* **496** 605
- [37] Bashinsky S and Bertshinger E 2002 *Phys. Rev. Lett.* **87** 1301
- [38] Springel V et al 2005 *Nature* **435** 629
- [39] Durrer R Gabrielli A Joyce M and Sylos Labini F 2003 *Astrophys.J.* **585** L1
- [40] Jiménez J B and Durrer R 2010 arXiv:1006.2343v1
- [41] Peacock J A 1999 “*Cosmological physics*” (Cambridge University Press, Cambridge)
- [42] Saslaw W C 2000 “*The Distribution of the Galaxies*” (Cambridge University Press)
- [43] Kerscher M 1999 *Astron.Astrophys.* **343** 333
- [44] Vasilyev N L Baryshev Yu V and Sylos Labini F 2006 *Astron.Astrophys.* **447** 431
- [45] Strauss M A et al 2002 *Astrophys.J.* **124** 1810
- [46] Eisenstein D J et al 2001 *Astrophys.J.* **12** 2267
- [47] Blanton M R and Roweis S 2007 *Astron.J.* **133** 734
- [48] Sylos Labini F 2010 arXiv:1011.4855v1
- [49] Madgwick D S et al 2002 *Mon.Not.R.Acad.Soc* **333** 133
- [50] Gabrielli A and Sylos Labini F 2001 *Europhys.Lett.* **54** 1
- [51] Kim J Park C Gott J R and Dubinski J 2009 *Astrophys.J.* **701** 1547
- [52] Eisenstein D J et al 2001 *Astronom.J.* **12** 2267
- [53] Fisher R A and Tippet L H C 1928 *Cambridge Phil. Soc.* **28** 180
- [54] Gumbel E J 1958 *Statistics of Extremes* (Columbia University Press)
- [55] Bramwell S T 2009 *Nature Physics* **5** 444
- [56] Landy S D and Szalay A 1993 *Astrophys.J.* **412** 64
- [57] Davis M and Peebles P J E 1983 *Astrophys.J.* **267** 465
- [58] Park C Vogeley M S Geller M J and Huchra J P 1994 *Astrophys.J.* **431** 569
- [59] Benoist C Maurogordato S da Costa L N Cappi A and Schaeffer R 1996 *Astrophys.J.* **472** 452
- [60] Zehavi I et al 2002 *Astrophys.J.* **571** 172
- [61] Sawangwit U et al 2009 arXiv:0912.0511v1
- [62] Martínez V J et al 2009 *Astrophys.J.* **696** L93
- [63] Bondi H 1952 *Cosmology* (Cambridge University Press, Cambridge)
- [64] Clifton T and Ferreira P G 2009 *Phys.Rev.* **D80** 103503
- [65] Sylos Labini F 2010 *Astron.Astrophys.* **523** A68
- [66] Lavaux G Tully R B Mohayaee R and Colombi S 2010 *Astrophys.J.* **709** 483
- [67] Watkins R Feldman H A and Hudson M J 2009 *Mon.Not.R.Acad.Soc.* 392 743
- [68] Kashlinsky A Atrio-Barandela F Kocevski D and Ebeling, H 2008 *Astrophys.J.* **686** L49
- [69] Kashlinsky A Atrio-Barandela F Ebeling H Edge A and Kocevski D. 2010 *Astrophys.J.* **712** L81
- [70] Ellis G 2008 *Nature* **452** 158
- [71] Buchert T 2008 *Gen.Rel.Grav.* **40** 467
- [72] Wiltshire D L 2007 *Phys.Rev.Lett.* **99** 251101
- [73] Clarkson C and Maartens R 2010 *Class. Quantum Grav.* **27** 124008
- [74] Rasanen S 2008 *Int.J.Mod.Phys* **D17** 2543
- [75] Kolb E W Marra V and Matarrese S 2010 *Gen.Rel.Grav.* **42** 1399
- [76] Célérier M N Bolejko K and Krasiński A 2010 *Astron. Astrophys.* **518** A21
- [77] Baryshev Y V 2006 *AIP Conf.Proc.* **822** 23
- [78] Joyce M et al 2000 *Europhys.Lett.* **50** 416